

Predictive Coding Systems for Electronic Discovery

Dhivya Soundarajan, HCI Master's degree student

Professor Sara Anne Hook, M.B.A., J.D.

HCC Department Brown Bag

October 14, 2016



WHAT IS ELECTRONIC DISCOVERY (E-DISCOVERY)?

- Electronic discovery (e-discovery) is something that impacts everyone, whether they know if or not, because it deals with the proper collection, preservation, analysis and production of evidence in digital form.
- To put it bluntly, if you are sued, the opposing party's lawyer will be requesting nearly every piece of digital evidence in any format that might be relevant to the case (including social media).
- This presentation will concentrate on the use of predictive coding in civil cases, but e-discovery is part of criminal cases as well as other types of audits and investigations. For example, see *Clark v. State*, 915 N.E.2d 126 (2009) – an early Indiana Supreme Court case that allowed the defendant's MySpace page to be admitted into evidence in spite of various objections by his lawyer.
- An especially important issue for anyone in the Informatics, Media Arts and IT industries.
- Anyone can find himself/herself needing to comply with requests for potentially relevant evidence – in electronic or paper form.

HISTORY OF ELECTRONIC DISCOVERY IN THE U.S.

- Series of decisions in *Zubulake v. UBS Warburg* and the 2006 amendments to the Federal Rules of Civil Procedure, a new area within law practice appeared, the law regarding electronic discovery (e-discovery).
- The phase of litigation known as discovery has existed for many years, with opposing parties and their lawyers making requests and exchanging documents that are relevant to a case.
- E-discovery transformed this process from the paper-based, pre-Internet world of discovery to a whole series of rules and decisions related to how to identify, collect, preserve, analyze, review, produce and present electronically-stored information (ESI).
- New e-discovery industry developed – and now there are vendors who offer systems and software that integrates e-discovery, digital forensics and litigation support systems.
- Efforts to determine standards and best practices, with EDRM being one example, along with the proclamations and guidelines issues by The Sedona Conference.
- Statistics indicate the career opportunities in e-discovery as well as information governance are going to significantly increase in the future.

ELECTRONIC DISCOVERY REFERENCE MODEL (EDRM)

Electronic Discovery Reference Model

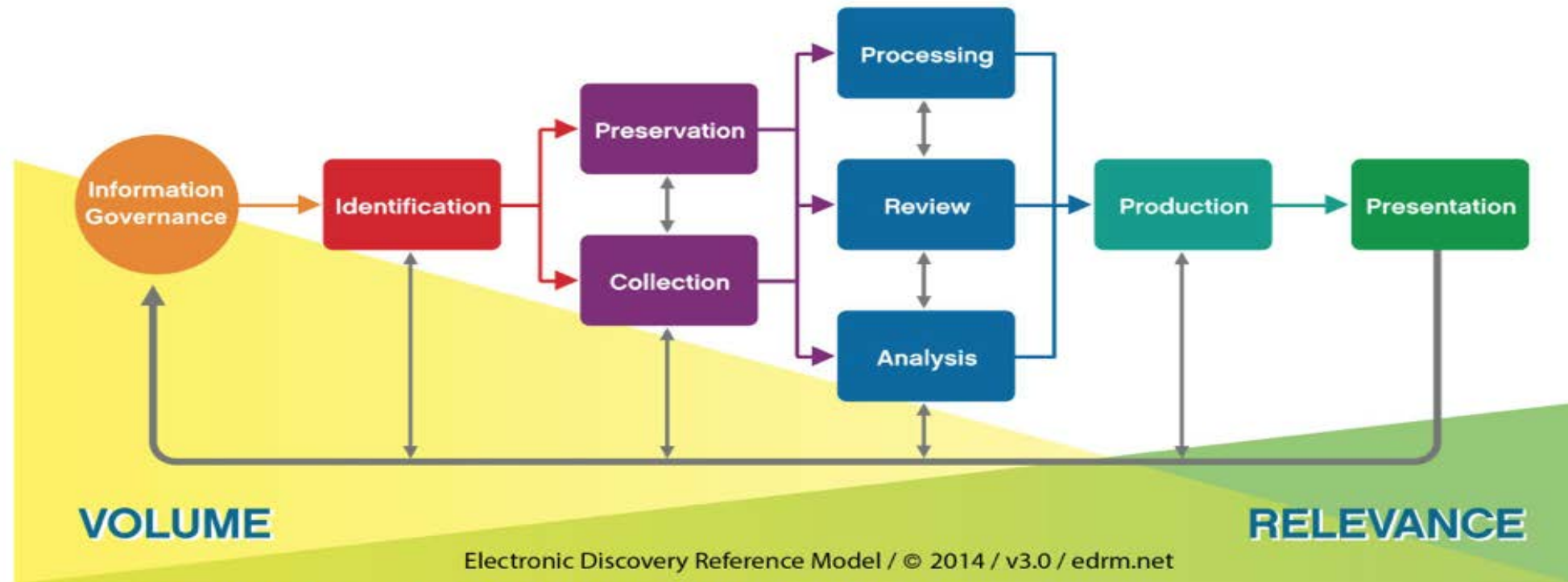


Figure 1 Diagram of the Electronic Discovery Reference Model

See <http://www.edrm.net/resources/edrm-stages-explained>, accessed 10/11/16.

E-DISCOVERY CHALLENGES

- Not only is this evidence now primarily in digital form, but it also exists a wide range of media and formats, from word processing and spreadsheet files to photographs, blog postings, videos, emails and websites.
- The terminology Electronically Stored Information (ESI) was chosen to reflect current and potential future technologies and cast a wide net in the discovery process.
- As noted in a recent survey conducted by Exterro, Inc., data volume is still the largest obstacle in e-discovery, with the second biggest obstacle being identifying and accessing sources of ESI.
- Recent debates and court decisions have focused on ESI that is posted on social media sites (Facebook) and text messages.
- Informal and transient communications beyond text messages, including new services for mobile devices and messaging apps, as well as data from wearable technology (fitness trackers) and the Internet of Things.
- The Federal Rules of Civil Procedure (FRCP), which govern courts in the federal court system, were revised again in December 2015, with an emphasis on proportionality, streamlining the process and cooperation with clarification of when and what types of sanctions can be imposed for spoliation and other intentional conduct.

WHAT IS PREDICTIVE CODING?

- Predictive coding is the use of keyword search, filtering and sampling to automate portions of an e-discovery process, especially the review stage.
- The goal of predictive coding is to reduce the number of irrelevant and non-responsive ESI that needs to be reviewed manually.
- May also be called – or part of – Technology-Assisted Review (TAR)
- A faulty and incomplete e-discovery process, particularly during the review stage, can result in sanctions and waive the attorney-client privilege or other confidentiality doctrine.
- Such failures, especially for breaches in confidentiality, can result in disciplinary action being taken against the lawyer by the state or states where he/she is licensed.
- Predictive coding systems can assist with the overall e-discovery process, leaving the humans to concentrate on reviewing the remaining set of ESI before it is produced to the opposing party.
- “[r]esearch shows that human review is far from perfect.” *Dynamo Holdings Ltd. P’ship v. Comm’r of Internal Revenue*, WL 4204067 (T.C. July 13, 2016).

COMMON TOOLS IN PREDCTIVE CODING?TAR

- Concept searching
- Contextual searching
- Metadata searching (ESI must usually be produced in native format with the metadata intact)
- Relevance probability and ranking
- Clustering
- Sorting ESI by issues

IS PREDICTIVE CODING ACCEPTED AS PART OF LITIGATION?

- Initially, predictive coding/TAR tools were looked at with considerable suspicion, even though information retrieval, indexing, machine learning and data analytics had been used in other disciplines for many years.
- The reticence to use these types of systems has faded, as illustrated by a long line of cases, starting with the strong support of computer-assisted review articulated in *Da Silva Moore v. Publicis Groupe*, described as the first published opinion recognizing TAR as “an acceptable way to search for relevant ESI in appropriate cases.”
- Summaries of recent cases about predictive coding/TAR can be found in The Sedona Conference’s new publication, *TAR Case Law Primer*.
- Cases indicate that judge’s will likely approve a party’s request to use predictive coding, absent some compelling objection.

HOW IS PREDICTIVE CODING USED IN LITIGATION?

- Early case assessment
- Reviewing client ESI before production
 - Prioritizing pre-production review
 - Sorting ESI by potential privilege
 - Quality control – comparing human review with predictive coding results
- Reviewing production from the opposing party
- Other stages of litigation, such as preparing for depositions, responding to summary judgment motions and working with expert witnesses

STATUS OF PREDICTIVE CODING

- “Overall, although the practice of predictive coding is still in its infancy, the number of courts addressing the issue is clearly on the rise. Courts seem to be moving towards permitting, but not requiring, this technology. Litigants that take reasonable positions and strive to work through their disputes with their opponents will typically be much better positioned to prevail in a predictive coding dispute.” (Wallis M. Hampton, Predictive Coding: It’s Here to Stay. *E-Discovery Bulletin*, June/July 2014, https://www.skadden.com/sites/default/files/publications/LIT_JuneJuly14_EDiscoveryBulletin.pdf, accessed 10/10/16.)
- The support for predictive coding has increased in the past two years since this article was published.

INTRODUCING DHIVYA SOUNDARAJAN

- Over the past year, Dhivya Soundarajan, a master's-level student in Human-Computer Interaction (HCI), has been designing a simple predictive coding system for me based on readily-available software and natural language processing.
- Dhivya will share the process of developing the system, the software she used and what she has designed so far.
- She will also discuss our future work, which will include usability testing of the system with a focus group of lawyers who are responsible for e-discovery and the features and functionality that we would like to add to the system.

FUNCTIONS OF PREDICTIVE SYSTEMS

Multimodal Input

Uses different types of unstructured text, digital archives, emails etc.

Concept Search

An automated information retrieval method to search electronically stored unstructured text which are conceptually similar to the information provided in a search query.

Supervised Machine Learning

The system should not only depend on passive analysis of data but should accept the lawyer's periodic input to enhance the system's efficiency.

FUNCTIONS OF PREDICTIVE SYSTEMS

Distributes Data

Data is divided into ranges and distributed to multiple servers.

Optimized Storage and Retrieval

Parallel processing of huge amount of data.

Ensures Easy Access of the System

Interface is very simple and more usable with different filter options, sort functionalities.

FEATURES OF PREDICTIVE SYSTEMS

Increases Accuracy

Overcomes the problem of
Boolean keyword searching.
Brings down false positives.
Starts including false negatives.

Provides Assurance

Since there is periodic
update/feedback from the
lawyer.

Eliminates Manual Filtering

Can make the documents
protected.
This is essential in preserving
attorney-client privilege and
attorney work-product.

FEATURES OF PREDICTIVE SYSTEMS

Provides Reliability

Load balancing of data - builds data easily in case of system failure.

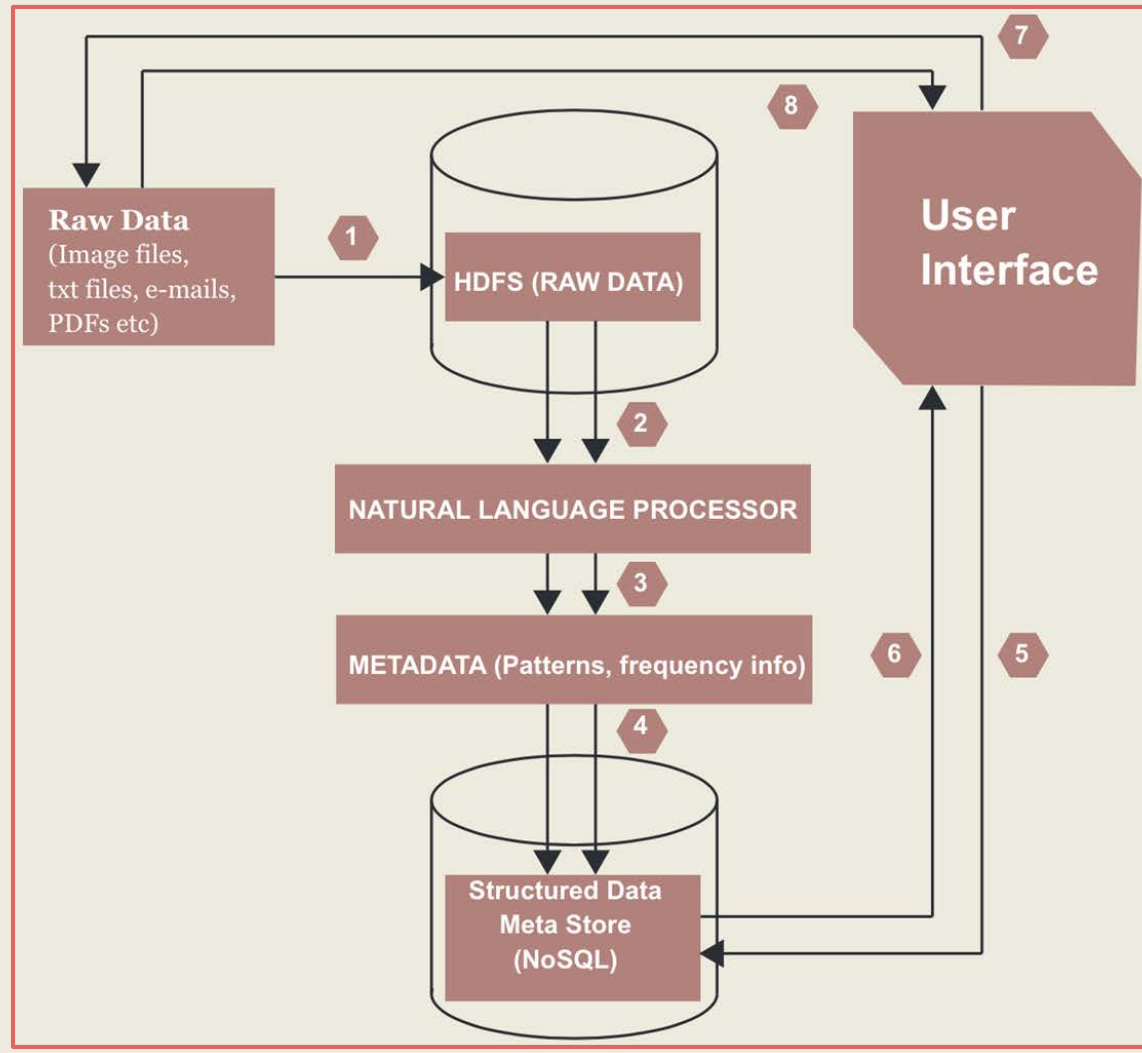
Robust

Huge set of unstructured data can be handled easily

SYSTEM ARCHITECTURE

System Architecture:

1. The raw data will be stored in the Hadoop file system along with its location mapping .
2. The data will be extracted for analyzing the patterns.
3. Metadata will be generated for each file based on the computational and processing algorithms.
4. The generated metadata will be stored as structured data in the meta store.
5. The user request will be passed to the meta store
6. The result will be the list containing the documents name and their source location.
7. The user passes the request along with the location parameters
8. The exact files are passed as output to the front-end of the system.



MACHINE LEARNING MODULE - NLP

TRAIN

ANALYSE

EVALUATE

Use several subsets of files (control sets) that are quintessential, identified by well-trained professionals for both the following cases in order to calibrate the system.

- Positive Sets - relevant files
- Negative Sets - irrelevant files

Then use training sets to train the system.

Apply the identified appropriate filters, classifiers and use the pre existing models with tailored specifications to analyse the system.

(ex.SMO models etc; Tools - Weka etc)

Check for

- Precision
- Recall
- F- Measure

Since the system is under a supervised learning, system training should happen periodically with new training sets as per the requirement. Then finalize the model for the system.

PRECISION

RECALL

F-MEASURE

~~—Data (true Positive)—~~
 $\text{Data}(\text{true Positive} + \text{false Positive})$

~~—Data (true Positive)—~~
 $\text{Data}(\text{true Positive} + \text{false Negative})$

$2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$

HADOOP DISTRIBUTED FILE SYSTEM MAP-REDUCE

COLLECT

Collect different feeds from different nodes (distributed in the cloud).

Example:

- **Documents**
- **Text messages**
- **Emails**

PROCESS

Process data as it flows.

- **Calculate**
- **Transform**
- **Process**
- **Augment**

VISUALIZE/QUERY

Display processed files as results of user Query.

USER INTERFACE

SMART PREDICTION

Bankruptcy

SEARCH

BY NAME☐

BY DATE☐

BY CASE TYPE☐

MOST FREQUENT☐

BY YEAR☐

MY STATE☐

REFRESH

DOCUMENTS

NOT PROTECTED/
PROTECTED

Not Protected ☒ Relevancy

File 1	<input type="checkbox"/>
File 2	<input type="checkbox"/>
File 3	<input type="checkbox"/>
File 4	<input type="checkbox"/>
File 5	<input type="checkbox"/>
File 6	<input type="checkbox"/>
File 7	<input type="checkbox"/>
File 8	<input type="checkbox"/>

DOWNLOAD

SCREEN SHOTS

ANALYSIS RESULTS-WEKA TOOL

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **FilteredClassifier** -F "weka.filters.unsupervised.attribute.StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -token

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds **10**
☐ Percentage split % 66
More options...

(Nom) @@class@@

Start Stop

Result list (right-click for options)

16:17:34 - meta.FilteredClassifier

Classifier output

```
=== Run information ===  
  
Scheme:      weka.classifiers.meta.FilteredClassifier -F "weka.filters.unsupervised.attribute.StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -  
Relation:    _Users_dhivyasivasankar_Desktop_Pc  
Instances:   18  
Attributes:  3  
             text  
             filename  
             @@class@@  
Test mode:   10-fold cross-validation  
  
=== Classifier model (full training set) ===  
  
FilteredClassifier using weka.classifiers.functions.SMO -C 1.0 -L 0.001 -P 1.0E-12 -N 0 -V -1 -W 1 -K "weka.classifiers.functions.supportVector.PolyKernel -E 1  
  
Filtered Header  
@relation '_Users_dhivyasivasankar_Desktop_Pc-weka.filters.unsupervised.attribute.StringToWordVector-R1,2-W1000-prune-rate-1.0-N0-stemmerweka.core.stemmers.Nul  
  
@attribute @@class@@ {bfiles,cfiles}  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile1 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile10 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile2 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile3 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile4 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile5 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile6 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile7 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile8 numeric  
@attribute /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile9 numeric  
@attribute Bank numeric  
@attribute Bankruptcy numeric  
@attribute Cassy numeric  
@attribute Crime numeric  
@attribute Defects numeric  
@attribute Financial numeric  
@attribute Lawyers numeric  
@attribute London numeric  
@attribute Manchester numeric  
@attribute Merchant numeric
```

Status

OK

Log x 0

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **FilteredClassifier** -F "weka.filters.unsupervised.attribute.StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -token

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds
☐ Percentage split %
 More options...

(Nom) @@class@@

Start Stop

Result list (right-click for options)

16:17:34 - meta.FilteredClassifier

Classifier output

```

@attribute house numeric
@attribute italy numeric
@attribute judgement numeric
@attribute juristication numeric
@attribute law numeric
@attribute lawyer numeric
@attribute lawyers numeric
@attribute litigation numeric
@attribute loss numeric
@attribute manage numeric
@attribute mansion numeric
@attribute minimal numeric
@attribute mistakes numeric
@attribute owner numeric
@attribute penalties numeric
@attribute plaintiff numeric
@attribute pool numeric
@attribute pound numeric
@attribute property numeric
@attribute rate numeric
@attribute respect numeric
@attribute seal numeric
@attribute sell numeric
@attribute shop numeric
@attribute state numeric
@attribute suicide numeric
@attribute suit numeric
@attribute summon numeric
@attribute swimming numeric
@attribute tackle numeric
@attribute takeover numeric
@attribute trends numeric
@attribute workers numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/Untitled-20 numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/cfile1 numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/cfile2 numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/cfile3 numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/cfile6 numeric
@attribute /Users/dhivyasisvasankar/Desktop/Pc/cfiles/cfile7 numeric
  
```

Status

OK Log x 0

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier

Choose **FilteredClassifier** -F "weka.filters.unsupervised.attribute.StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -token

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds
☐ Percentage split %
 More options...

(Nom) @@class@@

Start Stop

Result list (right-click for options)

16:17:34 - meta.FilteredClassifier

Classifier output

BinarySMO

Machine linear: showing attribute weights, not support vectors.

```

-0.1483 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile1
+ -0.0988 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile10
+ -0.1371 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile2
+ -0.059 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile3
+ -0.063 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile4
+ -0.0145 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile5
+ -0.027 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile8
+ -0.0495 * (normalized) /Users/dhivyasivasankar/Desktop/Pc/bfiles/bfile9
+ -0.059 * (normalized) Bank
+ -0.4601 * (normalized) Bankruptcy
+ -0.1483 * (normalized) Cassy
+ -0.0052 * (normalized) Crime
+ -0.059 * (normalized) Defects
+ -0.1483 * (normalized) Financial
+ -0.1483 * (normalized) Lawyers
+ -0.154 * (normalized) London
+ -0.1045 * (normalized) Manchester
+ -0.027 * (normalized) Merchant
+ -0.059 * (normalized) Plaintiff
+ -0.059 * (normalized) Statelaw
+ -0.027 * (normalized) account
+ -0.0988 * (normalized) attorney
+ -0.027 * (normalized) balance
+ -0.0824 * (normalized) bank
+ -0.1371 * (normalized) bankruptcy
+ -0.027 * (normalized) business
+ -0.0495 * (normalized) cars
+ 0.0012 * (normalized) case
+ -0.0988 * (normalized) close
+ 0.0328 * (normalized) court
+ -0.1045 * (normalized) cris
+ -0.027 * (normalized) crucial
+ 0 * (normalized) current
+ -0.027 * (normalized) currernt
  
```

Status

OK Log x 0

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **FilteredClassifier** -F "weka.filters.unsupervised.attribute.StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -token

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds
☐ Percentage split %
 More options...

(Nom) @@class@@

Start Stop

Result list (right-click for options)

16:17:34 - meta.FilteredClassifier

Classifier output

```

+ 0.0119 * (normalized) stolen
+ 0.0309 * (normalized) suicide
+ 0.0328 * (normalized) trouble
+ 0.0817 * (normalized) trust
+ 0.0691 * (normalized) virtual
+ 0.0534

Number of kernel evaluations: 171 (94.737% cached)

Time taken to build model: 0.02 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      17          94.4444 %
Incorrectly Classified Instances    1           5.5556 %
Kappa statistic                    0.8889
Mean absolute error                 0.0556
Root mean squared error             0.2357
Relative absolute error             11.1765 %
Root relative squared error         47.1502 %
Total Number of Instances          18

=== Detailed Accuracy By Class ===

               TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
               0.900    0.000    1.000     0.900   0.947     0.894    0.950    0.956    bfiles
               1.000    0.100    0.889     1.000   0.941     0.894    0.950    0.889    cfiles
Weighted Avg.   0.944    0.044    0.951     0.944   0.945     0.894    0.950    0.926


=== Confusion Matrix ===

 a b  <-- classified as
 9 1 | a = bfiles
 0 8 | b = cfiles

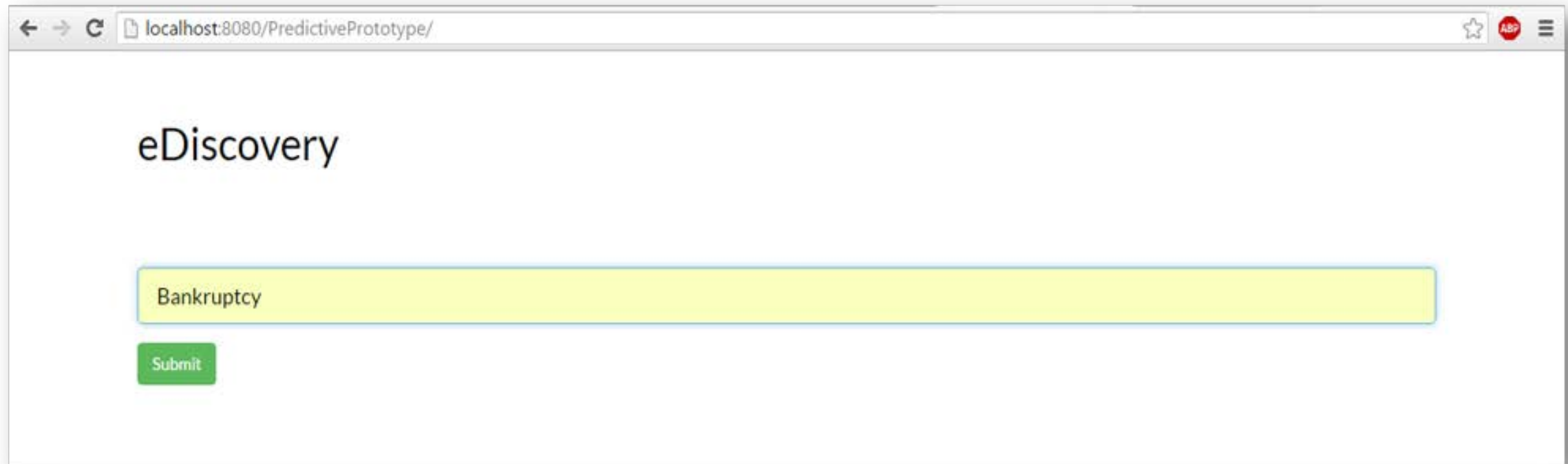
```

Status

OK

Log  x 0

SAMPLE WEB-BASED PROTOTYPE



← → ↻ localhost:8080/PredictivePrototype/ ☆ ASP ≡

eDiscovery

localhost:8080/PredictivePrototype/FileSearch

eDiscovery

Submit

Document Name	Download Link
AnnexAToSFAFinal	Download
AssignmentAndSecurityAgmtFinal52604	Download
DBOralOpinion	Download
DBOralOpinionSummary	Download
DTAFinal52804	Download
FPAS2604	Download
LitigationFacilityAgmtFinal52604	Download
PlanFinal52604	Download

FUTURE WORK

- As of now, we are working with an ideal set of data that we created.
- Now we have to gather some real data sets.
- Also work on integrating the overall modules – database, logic, Natural Language Processing.
- Test with a focus group of lawyers in the field of bankruptcy.
- Obtain data sets in other areas of the law.

Any Questions?

Thank you for attending the HCC Brown Bag today!

